



TITLE:

# 最適停止をともなう決定過程 (制御過程論 II)

AUTHOR(S):

古川, 長太

---

CITATION:

古川, 長太. 最適停止をともなう決定過程 (制御過程論 II). 数理解析研究所講究録 1970, 90: 171-207

ISSUE DATE:

1970-06

URL:

<http://hdl.handle.net/2433/108115>

RIGHT:

#### Ⅳ 決定過程

## 最適停止をともなう決定過程

九大理 古川長太

## §1. 序

離散時間確率的系に対する最適化問題では, stage 数を予め指定せずに無限 stage にわたっての政策を問題とするならば, 通常, 評価関数として無限和 (stage についての和), またはその 1 stage 当りの平均の型を採用する。ここでは無限和の型の評価関数を最大化することに限って考えることにする。

無限和の評価関数を採用することは, 一応, 我々が無限 stage にわたって行動し続けることを前提としているが, 現実にはいつかは stop せざるを得ない。また, おそらく stop した方がかえって有利な場合もある。また, いわゆる sequential analysis のように terminal decision が要求されるときは, 当然 いつかは stop しなくては問題の

解決にならない。

ここに一つの決定過程 ( $P$ ) が与えられたとし,  $P$  の状態空間を  $S$ , 行動空間を  $A$  とする。  $S$  に仮想的な状態  $s_0$

(吸収点) を加え,  $A$  に  $stop$  を表わす行動  $a_0$  を加えて,  $S' = S + \{s_0\}$ ,  $A' = A + \{a_0\}$  を新たにそれぞれ状態空間,

行動空間として新しい決定過程 ( $P'$ ) を構成すると,  $P'$  は  $stop$  するという行動をも含んだ決定過程になる。  $P'$  に対して無限和の評価関数を採用すると, 表面上は無限和であるが  $a_0$  をとった時点で process は実際には  $stop$  し, したがって, はじめに述べた問題は, このような定式化により通常の最適化問題に帰着されるかのように見える。しかしながら,  $P'$  における最適政策が  $P$  において果たして確率1で  $stop$  することになるかどうか分らないし, また  $P$  で  $stop$  する確率を計算することも一般に仲々困難である。

以上の理由により, ここに全く新しい定式化で考えなおすことにする。この報告では“確率1で  $stop$  する政策” (非常に rough な表現であるが) を考え, そのような政策の集合の中で最適なものを追究する。正確な定義は次の章にゆずるとして, さらに詳しく述べれば, 通常の意味の政策  $\pi$  と, 停止時間  $\tau$  を pair として考え, あらゆる pair  $(\pi, \tau)$  の集合の中で与えられた評価関数を最大化

することを問題とする。停止時間は停止規則と同等であるから、 $\mu$  の政策  $(\pi, \tau)$  はまた、通常の意味の政策と停止規則との  $\mu$  とみなせる。

近來、ソ連でマルコフ過程論(特に強マルコフ過程)の発展に伴って、E. B. Dynkin および A. N. Shiryaev 等がいくつかの具体的な最適停止問題について、興味ある結果を与え、またアメリカでは古くは Y. S. Chow, H. Robbins, 近年では L. E. Dubins, L. J. Savage 等により多くの具体的な最適停止問題が扱われて来た。

この報告は、定式化としては、前述のように通常の stochastic control problem に停止規則を入れた型でもあるし、また、別の面から見ると、上記の最適停止問題の一般化にさらに通常の control action を入れた型ともみなせる。

この報告は全体を通じて、上記の意味での  $\mu$  の政策の class における最適政策(最適性の規準は  $\delta$  でいくつか与えられる)の存在定理と、最適政策の性質について論じたものである。その他、これらに付随した多くの定性的な興味ある結果が与えられている。最適解の構成に関しては、状態空間、行動空間をごく簡単なものに限定すれば容易であるが、この報告でとり上げたような一般的な距離空間

では、可成り困難な問題で、この報告でも explicit な形で述べられ所までは追究されていない。(存在定理の証明方法が有限 stage による政策改良の線に沿っているから、この意味では可成り constructive である)

以上のように、最適政策の構成法に関しては不十分であるが、存在定理 および最適政策の性質等に関しては、一応初期の目的を達したので、この報告で総括する。

## §2. 記号および定式化

### 2.1 一般的な定義

Borel set ; complete separable metric space の Borel subset  
( $X, Y$  等で表わす)

Borel set  $X$  上の probability measure ;  $X$  の Borel field 上の  
probability measure

$P(X)$  ;  $X$  上のすべての probability measure からなる空間

$\mathcal{G}(Y|X)$  ;  $x$  を与えたときの  $Y$  上の conditional prob. measure で

(i) 各  $x \in X$  に対し  $\mathcal{G}(\cdot|x)$  は  $Y$  上の prob. measure

(ii) 各 Borel subset  $B \subset Y$  に対し  $\mathcal{G}(B|\cdot)$  は  $X$  上の Baire 関数

$Q(Y|X)$ ; 定義空間  $X, Y$  をきめて 上のようなすべての  
conditional prob. measure からなる空間

$M(X)$ ; (i) §2 ~ §4 では  $X$  上のすべての実数値 Baire  
関数からなる空間

(ii) §5, §6 では  $X$  上のすべての有界 Baire 関数  
からなる空間

$u, v \in M(X)$  に対して  $u \geq v$ ;  $u(v) \geq v(x)$  for all  $x \in X$

$p \in P(X), u \in M(X)$  に対して  $pu$ ;  $pu = \int_X u(v) dp(x)$

$g \in Q(Y|X), u \in M(XY)$  に対して  $gu$ ;  $gu(v) = \int_Y u(x, y) dg(y|x)$

$p \in P(X)$  が degenerate;  $p\{x_0\} = 1$  for some  $x_0 \in X$

$g \in Q(Y|X)$  が degenerate;  $g(\cdot|x)$  が  $x$  に対して degenerate

$g \in Q(Y|X)$  が degenerate であれば  $g(\{f(x)\}|x) = 1$  for

all  $x \in X$  なる  $X$  から  $Y$  への Baire 関数  $f$  が存在する

ことになる, 従って  $u \in M(XY)$  なる  $u$  に対しては

$$fu(v) = u(x, f(v)) \quad \text{for all } x \in X$$

となる。

## 2.2 決定過程に固有な定義

停止規則をとるような決定過程に対する最適化問題は  
つぎの 6 箇の要素  $S, A, g, r, q, \alpha$  で決定される。

$S$  ; 状態空間 (state space) を表わし, ある空でない Borel set

$A$  ; 行動空間 (action space) を表わし, ある空でない Borel set

$\mathcal{F} = \{\mathcal{F}_1, \mathcal{F}_2, \dots\}$  ; system の運動法則を表わし, 各  $n$  につき

$$\mathcal{F}_n \in \mathcal{Q}(S | H_n A), \quad H_n = S A S A \dots S \quad (2n-1 \text{ 回})$$

$r \in M(SA)$  ; 利得関数 (reward function)

$g$  ; 最終利得関数で  $S$  上の有界 Baire 関数

$\alpha$  ; 割引率 (discount factor)

$$(i) \quad \S 2 \sim \S 4 \text{ では } 0 \leq \alpha \leq 1$$

$$(ii) \quad \S 5, \S 6 \text{ では } 0 \leq \alpha < 1$$

つぎに, policy, stopping time および stopped policy 等についての定義を述べる.

政策 (policy)  $\pi = \{\pi_1, \pi_2, \dots\}$  ;  $\pi_n \in \mathcal{Q}(A | H_n)$   $n=1, 2, \dots$

Markov policy ;  $\pi = \{\pi_1, \pi_2, \dots\}$  において各  $\pi_n$  が  $\mathcal{Q}(A | S)$  の degenerate な要素

Markov policy を  $\{f_1, f_2, \dots\}$  ( $k \geq 1$   $f_k$  は  $S$  から  $A$  への Baire 関数) と書く.

定常政策 (stationary policy) ;  $\pi = \{f_1, f_2, \dots\}$  が

Markov policy かつ  $f_n = f$  for all  $n$

stationary policy を  $f^\infty$  と書く.

$\pi = \{\pi_1, \pi_2, \dots\}$  に対して  $n\pi = \{\pi_{n+1}, \pi_{n+2}, \dots\}$  と書く. 故に  $0\pi = \pi$



$\pi$  を任意の policy とすれば,  $\pi$  と  $g$  とで "つぎ" のような conditional probability measure  $e_\pi$  が定義される。

$$e_\pi = \pi_1 g_1 \pi_2 g_2 \cdots \in Q(SASAS \cdots | S)$$

$\mathcal{G}$  ;  $S$  の Borel field

$\mathcal{A}$  ;  $A$  の Borel field

標本空間 (sample space) ;  $\Omega = SASAS \cdots$

標本点 (sample point) 又は 道 (path) ;  $\omega = (s_1, a_1, s_2, a_2, \dots)$

$\mathcal{F}_n$  ;  $\{\omega \mid s_1 \in E_1, a_1 \in F_1, s_2 \in E_2, a_2 \in F_2, \dots, s_n \in E_n\}$

(ただし  $E_i \in \mathcal{G}, F_i \in \mathcal{A}, i=1, 2, \dots, n$ )

の型の  $\omega$ -集合のあらゆる和の族

( $\mathcal{F}_n$  は 明らかに  $\sigma$ -field になる)

policy  $\pi$  に付随する停止時間 (stopping time)

$$t = t(\omega) ; \begin{cases} \text{(i) 正整数値の確率変数} \\ \text{(ii) } e_\pi[\{t(\omega) < \infty\} | s_1] = 1 \text{ for all } s_1 \in S \\ \text{(iii) } \{t(\omega) = n\} \in \mathcal{F}_n \text{ for each } n \end{cases}$$

$C(\pi)$  ;  $\pi$  に付随するすべての stopping time の集合

停止政策 (stopped-policy)  $(\pi, t)$  ;

$\pi$  は任意の policy,  $t \in C(\pi)$  のとき  $(\pi, t)$  を stopped-policy といい, 簡単のために  $s$ -policy とかく

$C^N(\pi)$  ;  $C(\pi)$  に属しかつ

$$E_{\pi}[\{t(\omega) \leq N\} \mid s_1] = 1 \quad \text{for all } s_1 \in S$$

とみたすすべての stopping time の集合

打ち切政策 (truncated-policy) ;  $s$ -policy  $(\pi, t)$  として

あってかつ  $t \in C^N(\pi)$  なる正整数  $N$  が存在するものを

を truncated-policy といい、簡単のために  $t$ -policy

とかく

$$\mathcal{B}_n; \{ \omega \mid s_1 \in E_1, s_2 \in E_2, \dots, s_n \in E_n \} \quad (\text{ただし } E_i \in \mathcal{G} \quad i=1, 2, \dots, n)$$

の型の  $\omega$ -集合のあらゆる和の族

( $\mathcal{B}_n$  は明らかに  $\sigma$ -field)

$$\mathcal{G}_n; \{ \omega \mid s_n \in E \} \quad (\text{ただし } E \in \mathcal{G})$$

の型のあらゆる  $\omega$ -集合の族

$t \in C(\pi)$  が Markov stopping time ;

$$\left\{ \begin{array}{l} \text{(i)} \quad \{t(\omega) = n\} \in \mathcal{B}_n \quad \text{for each } n \\ \text{(ii)} \quad \{t(\omega) = n\} = \{t(\omega) > n-1\} \cap \Delta_n \end{array} \right.$$

$$\text{ただし } \Delta_n \in \mathcal{G}_n, \quad \text{for each } n$$

Markov  $s$ -policy  $(\pi, t)$  ;  $\left\{ \begin{array}{l} \text{(i)} \quad \pi \text{ が Markov policy} \\ \text{(ii)} \quad t \in C(\pi) \text{ かつ } t \text{ が Markov stopping time} \end{array} \right.$

stationary  $t$ -policy  $(\pi, t)$  ;  $(\pi, t)$  が Markov  $s$ -policy としてかつ,  $\pi$  が stationary policy

$t \in C(\pi)$  が stationary stopping time ;  $t$  が Markov で

かつ  $\Delta_n$  が  $n$  に無関数

つぎに評価関数の設定と, いくつかの最適性の規準を与える。

$$X_n = \sum_{k=1}^{n-1} \alpha^{k-1} r_{\Delta_k, a_k, \Delta_{k+1}} + \alpha^{n-1} g(\Delta_n), \quad n=1, 2, \dots$$

これを簡単のために

$$X_n = \sum_{k=1}^{n-1} \alpha^{k-1} r_k + \alpha^{n-1} g_n, \quad n=1, 2, \dots$$

とかくこともある。

$E^\pi$  ;  $e_\pi$  に関する expectation を表わす積分作用素

$\delta$ -policy  $(\pi, t)$  から得られる総期待利得を評価関数として採用する。すなわち

$$\begin{aligned} E^\pi(X_t) &= E^\pi \left[ \sum_{k=1}^{t-1} \alpha^{k-1} r_k + \alpha^{t-1} g_t \right] \\ &= e_\pi \left[ \sum_{k=1}^{t-1} \alpha^{k-1} r_k + \alpha^{t-1} g_t \right] \end{aligned}$$

我々の目的は, 上に定義した  $E^\pi(X_t)$  を, 何らかの意味で  $(\pi, t)$  について最大化することである。

$\Pi^N$  ;  $N$  番 stage までに対するすべての policy  $\{\pi_1, \pi_2, \dots, \pi_N\}$  の集合

$$\Lambda^N \equiv \{(\pi, t) \mid \pi \in \Pi^{N-1}, t \in C^N(\pi)\}$$

$\Pi$  ; 無限 stage にかたつてのすべての policy の集合

$$\Lambda \equiv \{(\pi, t) \mid \pi \in \Pi, t \in C(\pi)\}$$

$(p, \varepsilon)^N$ -optimal ;  $(\hat{\pi}, \hat{k}) \in \Lambda^N$  かつ

$$p\{E^{\hat{\pi}}(x_{\hat{k}}) \geq E^{\pi}(x_t) - \varepsilon\} = 1 \text{ for } \forall (\pi, t) \in \Lambda^N$$

のとき  $(\hat{\pi}, \hat{k})$  を  $(p, \varepsilon)^N$ -optimal といい

$(p, \varepsilon, \delta)$ -optimal ;  $(\hat{\pi}, \hat{k}) \in \Lambda$  かつ

$$p\{E^{\hat{\pi}}(x_{\hat{k}}) \geq E^{\pi}(x_t) - \varepsilon\} \geq 1 - \delta \text{ for } \forall (\pi, t) \in \Lambda$$

のとき  $(\hat{\pi}, \hat{k})$  を  $(p, \varepsilon, \delta)$ -optimal といい

$(p, \varepsilon)$ -optimal ;  $(\hat{\pi}, \hat{k})$  が  $(p, \varepsilon, 0)$ -optimal のとき, これ

を  $(p, \varepsilon)$ -optimal といい

$\varepsilon$ -optimal ;  $(\hat{\pi}, \hat{k}) \in \Lambda$  かつ

$$E^{\hat{\pi}}(x_{\hat{k}}) \geq E^{\pi}(x_t) - \varepsilon \text{ for } \forall (\pi, t) \in \Lambda$$

のとき  $(\hat{\pi}, \hat{k})$  を  $\varepsilon$ -optimal といい

optimal ;  $(\hat{\pi}, \hat{k})$  が 0-optimal のとき これを optimal

といい

$(p, \varepsilon)$ -dominate ;  $p\{E^{\hat{\pi}}(x_{\hat{k}}) \geq E^{\tilde{\pi}}(x_{\tilde{k}}) - \varepsilon\} = 1$  のとき

$(\hat{\pi}, \hat{k})$  は  $(\tilde{\pi}, \tilde{k})$  を  $(p, \varepsilon)$ -dominate する

といふ

§ 3.  $(p, \varepsilon, \delta)$ -optimal policy および  $(p, \varepsilon)$ -optimal policy の  
存在

先づ本論に入る前に  $\{x_n\}$ -process から  $\{\beta_n^N(\pi)\}$ -process

を構成し それにもとづく2つの proposition をあげる。

$\pi$  を任意の policy として各  $N \geq 1$  に対し,  $\{\beta_n^N(\pi)\}$ -process を次式で backward に定義する。

$$\begin{cases} \beta_n^N(\pi) = \max [x_n, E^{n-1, \pi}(\beta_{n+1}^N(\pi))] & n=1, 2, \dots, N-1 \\ \beta_N^N(\pi) = x_N \end{cases} \quad (3.1)$$

ただし 各  $x_n$  は  $E^\pi$  に関して可積分であるとし,  $E^{n-1, \pi}$  は conditional prob. measure  $\pi_n \otimes_n \pi_{n+1} \otimes_{n+1} \dots$  に関する積分を表わす。したがって policy に関する仮定から, 各  $N$ ,  $n$  につき  $\beta_n^N(\pi)$  は  $\mathcal{F}_n$ -可測である。

$\{\beta_n^N(\pi)\}$  を用いて  $\tau_N$  を定義する。

$$\tau_N(\pi) = \text{the first } n \text{ such that } \beta_n^N(\pi) = x_n. \quad (3.2)$$

(3.1) により  $\beta_N^N(\pi) = x_N$  だから  $\tau_N(\pi) \in C^N(\pi)$ 。

今后簡単のために, 誤解の恐れのない限り  $\tau_N(\pi)$  を  $\tau_N$  で表わす。  $\tau_N$  の定義から明らかに

$$x_{\tau_N} = \beta_{\tau_N}^N(\pi).$$

Proposition 1 (3.1), (3.2) で定義された  $\{\beta_n^N(\pi)\}$ ,  $\tau_N$  に対して

$$(a) \quad E^\pi(x_{\tau_N}) = \sup_{t \in C^N(\pi)} E^\pi(x_t)$$

$$(b) \quad E^\pi(\beta_{\tau_N}^N(\pi)) = \sup_{t \in C^N(\pi)} E^\pi(\beta_t^N(\pi))$$

Proposition 2 (3.1), (3.2) で定義された  $\{\beta_n^N(\pi)\}$ ,  $\tau_N$  に対し

$$E^\pi(x_{\tau_N}) = E^\pi[\beta_{\tau_N}^N(\pi)] = \beta_1^N(\pi)$$

(注意) Proposition 1 および Proposition 2 は, この報告の本束の目的のためには全く準備的な事に過ぎないが, これ自体また別の意味で非常に興味ある結果である。これらの証明はかなりの補助定理を必要とするのでここでは省略するが,  $\{\beta_n^N(\pi)\}$ -process が super Martingale 性をもつことから martingale system theory を用いることを付記しておく。

$$X^N \equiv ASAS \cdots S \quad (2N \text{ factors})$$

$\Sigma^*$ ;  $P(X^N)$  の上の  $\sigma$ -field で: 各 Borel subset  $B \subset X^N$  に

対し, 写像  $\nu(B); P(X^N) \rightarrow R^1$  が  $\Sigma^*$ -可測となる

最小の  $\sigma$ -field

$\Sigma^*$  は, weak-topology に関する  $P(X^N)$  の Borel field と一致することがある。

$$\Gamma^N \equiv \{(\delta, \nu) \mid \delta \in S, \nu = e_\pi(\delta) \text{ for some } \pi \in \Pi^N\}$$

Lemma 3.1  $\Gamma^N$  は  $S \cdot P(X^N)$  の Borel subset である。

次に  $\{\beta_n^N(\pi)\}$ -process に類似な方法で  $\{\beta_n^N(\nu)\}$ -process を構成する。

$\nu \in \mathcal{P}(X^N)$  を任意に与え,  $\nu$  の分解を考える.

$$\nu = \nu_1 \nu_2 \cdots \nu_{2N},$$

ただし

$$\nu_1 \in \mathcal{P}(A)$$

$$\nu_2 \in \mathcal{Q}(S|A)$$

$$\nu_{2n+1} \in \mathcal{Q}(A|ASAS \cdots AS) \quad (2n \text{ factors})$$

$$\nu_{2n} \in \mathcal{Q}(S|ASAS \cdots A) \quad (2n-1 \text{ factors}).$$

$\{\beta_n^N(\nu)\}$  を次式で backward に構成する.

$$\begin{cases} \beta_n^N(\nu) = \max[\alpha_n, \nu_{2n-1} \nu_{2n} \beta_{n+1}^N(\nu)], & n=1, 2, \dots, N-1 \\ \beta_N^N(\nu) = \alpha_N \end{cases} \quad (3.3)$$

ここで次の仮定をおく.

(A1) 上の (a), (b), (c) のうち少なくとも一つが成立するとは.

$$(a) \quad -\infty < \beta_1^N(\nu) < \infty \quad \text{for } \forall \nu \in \mathcal{P}(X^N)$$

$$(b) \quad \gamma \geq 0, \quad q \geq 0$$

$$(c) \quad \gamma \leq 0, \quad q \leq 0.$$

(注意) (A1) の (b), (c) はそれぞれ, いわゆる positive case, negative case である. また (a) は §5, §6 で扱う discounted case では常に満足される. discounted case を今后, D-case と呼ぶ.

$$\nu^{N*} \equiv \sup_{\pi \in \Pi^N} \beta_1^N(\pi).$$

一般に  $v^{N*}$  は可測 になるとは限らないが、次の Lemma が成立つ。

Lemma 3.2 (A1) を仮定すると、 $v^{N*}$  は絶対可測 である。

(略証)

$$v(\Delta, \nu) \equiv \beta_1^N(\nu)(\Delta)$$

$$B_\lambda \equiv \Gamma^N \cap \{(\Delta, \nu) \mid v(\Delta, \nu) > \lambda\}$$

とおくと  $\Sigma^*$  の作り方から  $v(\Delta, \nu)$  が  $SP(X^N)$  上で  $(\Delta, \nu)$  に関して可測 なることと、Lemma 3.1 により  $B_\lambda$  は  $SP(X^N)$  の Borel subset になる。一方

$$v^{N*}(\Delta) = \sup_{\nu \in \Gamma_\Delta^N} v(\Delta, \nu) \quad (\text{ただし } \Gamma_\Delta^N \text{ は } \Gamma^N \text{ の } \Delta\text{-section})$$

だから、 $C_\lambda = \{\Delta \mid v^{N*}(\Delta) > \lambda\}$  は  $B_\lambda$  の  $S \wedge$  の projection である。従って  $C_\lambda$  は analytic set。故に絶対可測 である。 (Q.E.D.)

Lemma 3.3 (A1) を仮定すると、任意の  $p \in P(S)$ 、任意の  $\varepsilon > 0$  に対して

$$p\{\beta_1^N(\hat{\pi}) \geq v^{N*} - \varepsilon\} = 1$$

なるような  $\hat{\pi} \in \Pi^N$  が存在する。

(略証)

Lemma 3.2 により任意の  $p \in P(S)$  に対して Borel set  $N_1 \subset S$  と可測関数  $v_0$  が存在して



$$p(N_1) = 0 \quad \text{かつ} \quad v_0(\omega) = v^{N^*}(\omega) \quad \text{for } \omega \notin N_1.$$

つぎに

$$\Gamma_\varepsilon \equiv \Gamma^N \cap \left[ \{(\omega, \nu) \mid \omega \notin N_1, \nu(\omega, \nu) > v_0(\omega) - \varepsilon\} \cup \{(\omega, \nu) \mid \omega \in N_1\} \right]$$

とおくと,  $\Sigma^*$  の定義と Lemma 3.1 により,  $\Gamma_\varepsilon$  は  $SP(X^N)$  の Borel subset である。また

$$\Gamma_\varepsilon \text{ の } \omega\text{-section} = \emptyset \quad \text{for all } \omega \in S$$

は,  $\Gamma^N$  と  $v^{N^*}$  の定義より明らか。故に G.W. Mackey の可測陰関数の定理 ([5] の Theorem 6.3) により, Borel set  $N_2 \subset S$  と可測関数

$$\gamma = \gamma(\omega) ; \quad S \rightarrow P(X^N)$$

が存在して,  $p(N_2) = 0$  かつ  $(\omega, \gamma(\omega)) \in \Gamma_\varepsilon$  for  $\omega \notin N_2$ .

$\Sigma^*$  の定義と  $\gamma$  の可測性から, 各  $B \subset X^N$  に対して

$$\gamma(\omega)(B) \text{ は } \omega \text{ に関して可測。故に } \mu(\cdot | \omega) = \gamma(\omega)(\cdot)$$

とおけば,  $\mu \in Q(X^N | S)$  となる。

$$\mu = \mu_1 \mu_2 \cdots \mu_{2N},$$

$$\mu_1 \in Q(A | S)$$

$$\mu_{2n} \in Q(S | SA \cdots A) \quad (2n \text{ factors})$$

$$\mu_{2n+1} \in Q(A | SA \cdots S) \quad (2n+1 \text{ factors})$$

と分解出来る。かつ

$$\mu_{2n}(\cdot | s_1 a_1 \cdots s_n a_n) = f_n(\cdot | s_1 a_1 \cdots s_n a_n) \quad \text{for } s_i \notin N_2$$

となる。  $\hat{\pi} = \{\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_N\}$  を

$$\hat{\pi}_n = \begin{cases} \mu_{2n-1} & \text{for } s_1 \notin N_2 \\ \pi'_n & \text{for } s_1 \in N_2 \end{cases} \quad \text{ただし } \pi'_n \text{ は任意} \\ (n=1, 2, \dots, N)$$

で定義すれば,  $s_1 \notin N_1 \cup N_2$  ならば

$$v(s_1, \delta(s_1)) = \beta_1^N(\hat{\pi})_{s_1} \geq v_0(s_1) - \varepsilon = v^{N*}(s_1) - \varepsilon.$$

・ しかるに  $p(N_1 \cup N_2) = 0$  だから

$$p\{\beta_1^N(\hat{\pi}) \geq v^{N*} - \varepsilon\} = 1 \quad (\text{Q.E.D.})$$

(注意) Lemma 3.3 の証明で Mackey の 閉関数の定理を用いたが, これが適用出来るためには  $P(X^N)$  が位相空間として standard Borel space であれば十分である。(S はもともと standard Borel space である)  $P(X^N)$  がこの § のはじめに導入した weak-topology に関して standard Borel space になることは K. R. Parthasarathy によって [6] の第 2 章に詳しく述べられている。

つぎの定理は 有限 stage の model に対する最適政策の存在定理であるが, 実はこれがこの § で求める  $(p, \varepsilon, \delta)$ -optimal  $\delta$ -policy の存在定理のための一つの重要な base である。

Theorem 3.1 (A1) を仮定する。任意の  $p \in P(S)$ , 任意の  $\varepsilon > 0$  に対して  $(p, \varepsilon)^N$ -optimal  $t$ -policy が存在する。

(略証)

Lemma 3.3 により 任意の  $p \in P(S)$ ,  $\varepsilon > 0$  に対して

$$p \{ \beta_1^N(\hat{\pi}) \geq v^{N*} - \varepsilon \} = 1$$

なる  $\hat{\pi} \in \Pi^{N-1}$  がある。故に

$$p \{ \beta_1^N(\hat{\pi}) \geq \beta_1^N(\pi) - \varepsilon \} = 1 \quad \text{for all } \pi \in \Pi^{N-1} \quad (3.4)$$

$\pi$  を任意の policy として,  $\tau_N(\hat{\pi})$ ,  $\tau_N(\pi)$  をそれぞれ (3.2) で定義される  $\hat{\pi}$  と  $\pi$  に対応する stopping time とする。

Proposition 2 により, (3.4) に代入して

$$p \{ E^{\hat{\pi}}(x_{\tau_N(\hat{\pi})}) \geq E^{\pi}(x_{\tau_N(\pi)}) - \varepsilon \} = 1 \quad \text{for } \forall \pi \in \Pi^{N-1}$$

Proposition 1-(a) により

$$p \{ E^{\hat{\pi}}(x_{\tau_N(\hat{\pi})}) \geq E^{\pi}(x_t) - \varepsilon \} = 1 \quad \text{for } \forall (\pi, t) \in \Lambda^N$$

(Q.E.D.)

ここで仮定 (A1) に代えて, 仮定 (B1) をおく。

(B1) 次の (a), (b), (c) のうち少なくとも一つが成立つこと。

(a) (A1) の (a) が各  $N \geq 1$  について成立つ

(b)  $\gamma \geq 0$ ,  $g \geq 0$

(c)  $\gamma \leq 0$ ,  $g \leq 0$

(注意) (B1) の (a) は D-case では常にみたされる。

Lemma 3.4 (B1) を仮定する。  $p \in P(S)$  の任意の要素,  $\varepsilon$  を任意の正数とする。各  $N \geq 1$  に対して  $(p, \varepsilon)^N$ -optimal  $\pi$ -policy を  $(\hat{\pi}_N, \hat{c}_N)$  とする。(その存在は Theorem 3.1 により保証される)。このとき一般性を失うことなく,  $p$  確率 1 で  $\{E^{\hat{\pi}_N}(x_{\hat{c}_N})\}$  は単調非減少であるとしてよい。

ここであらたに 次の仮定をおく。

$$(A2) \quad \sup_{(\pi, t) \in \Lambda} E^{\pi}(x_t) < \infty$$

$$(A3) \quad Y_k = Y'_k - Y''_k \quad (k=1, 2, \dots)$$

ただし  $Y''_k \geq 0$  で  $Y'_k, Y''_k$  は共に SAS 上の実数値

Baire 関数。

$$(A4) \quad E^{\pi}(x'_t) < \infty \text{ for every } (\pi, t) \in \Lambda$$

$$\text{ただし} \quad x'_t = \sum_{k=1}^{t-1} \alpha^{k-1} Y'_k + \alpha^{t-1} g_t.$$

$$(A5) \quad \{(x'_n)^-, n \geq 1\} \text{ が各 } c_{\pi} \text{ に関して一様可積分}$$

$$\text{ただし} \quad x'_n = \sum_{k=1}^{n-1} \alpha^{k-1} Y'_k + \alpha^{n-1} g_n.$$

(注意) (A2) は D-case では常にみたされる。また

(A3) において  $Y''_k \geq 0$  ということは  $Y''_k$  が cost または cost を含む損失の部分を表わすと思えばよい。このとき  $Y'_k$  には可測性の他に制限がないから, 我々の

問題は  $\text{cost}$  だけを評価関数として  $\text{cost}$  を最小にする最適化問題をも含むことになる。この意味で仮定 (A3) はごく一般的な仮定である。

以上の仮定のもとで次の Lemma を得る。

Lemma 3.5 (B1), (A2) ~ (A5) を仮定する。各  $N$  に対する  $(p, \varepsilon)^N$ -optimal  $t$ -policy を  $(\hat{\pi}^N, \hat{t}_N)$  で表わせば

$$p \left\{ \lim_{N \rightarrow \infty} E^{\hat{\pi}^N}(x_{\hat{t}_N}) \geq E^{\pi}(x_t) - \varepsilon \right\} = 1 \quad \text{for } \forall (\pi, t) \in \Lambda$$
 が成立する。

(畧証)

一般性を失うことなく,  $E^{\pi}(x_t) \neq -\infty$  としてよい。

(A4) より  $-\infty < E^{\pi}(x_t) \leq E^{\pi}(x'_t) < \infty$

$$\therefore -\infty < E^{\pi}(x''_t) < \infty \quad \text{ただし} \quad x''_t = \sum_{k=1}^{t-1} \alpha^{k-1} \gamma_k'' \quad (3.5)$$

一方  $t \in C(\pi)$  に対して

$$t > n \Rightarrow x_n \geq -(x'_n)^- - x''_t$$

故に

$$\int_{\{t > n\}} x_n^- d\pi \leq \int_{\{t > n\}} [(x'_n)^- + x''_t] d\pi \quad (3.6)$$

(3.5) (3.6) (A5) により

$$\lim_{n \rightarrow \infty} \int_{\{t > n\}} x_n^- d\pi = 0 \quad (3.7)$$

次に,  $t_N = \min(t, N)$  とおけば,  $(\pi, t_N) \in \Lambda^N$

次に  $(\hat{\pi}_N, \hat{c}_N)$  の定義により

$$\begin{aligned} \int_{\{t \leq N\}} x_t d\pi &= E^{\pi}(x_{t_N}) - \int_{\{t > N\}} x_N d\pi \\ &\leq E^{\hat{\pi}_N}(x_{\hat{c}_N}) + \varepsilon - \int_{\{t > N\}} x_N d\pi \quad \text{w.p.1} \\ &\leq E^{\hat{\pi}_N}(x_{\hat{c}_N}) + \varepsilon + \int_{\{t > N\}} x_N^- d\pi \quad \text{w.p.1} \end{aligned} \quad (3.8)$$

しかも (3.8) は各  $N \geq 1$  について成立つ。(3.8)において  $N \rightarrow +\infty$  とすれば (A2), Lemma 3.4, (3.7) により  $\lim_{N \rightarrow \infty} E^{\hat{\pi}_N}(x_{\hat{c}_N})$  は存在して finite で, Lemma の求める結果が得られる。 (Q.E.D.)

この節の主目的は  $(p, \varepsilon, \delta)$ -optimal  $s$ -policy の存在定理としてつきが与えられる。

Theorem 3.2 (B1), (A2) ~ (A5) を仮定する。

各  $p \in \mathcal{P}(S)$ , 各  $\varepsilon > 0$ , 各  $\delta > 0$  に対して  $(p, \varepsilon, \delta)$ -optimal  $s$ -policy が存在し, かつそれは  $\pi$ -policy の class 中にある。

(略証)

Lemma 3.5 に Egoroff の定理を適用すれば直ちに得られる。

Corollary 3.1      Theorem 3.2 と同じ仮定をおく。

(a)  $S$  が有限集合ならば, 任意の  $p \in P(S)$ , 任意の  $\varepsilon > 0$  に対して  $(p, \varepsilon)$ -optimal  $t$ -policy が存在する。

(b)  $p \in P(S)$  が finite support をもてば, 任意の  $\varepsilon > 0$  に対して  $(p, \varepsilon)$ -optimal  $t$ -policy が存在する。

(略証)

(a) または (b) の仮定があれば,  $p$ -確率 1 で  $\{E^{\pi_N}(x_{t_N})\}$  が一様収束することから明らかである。

#### §4. Markov stopped-policy

この § を通じて system の運動法則は Markov type であるとする。即ち

$$p_n \in Q(S|SA) \quad \text{for each } n$$

とする。ただし  $Q(S|SA)$  における  $|$  の右側の  $SA$  は  $n$  番 stage における状態空間と行動空間を表わすものとする。

§3 で  $(p, \varepsilon, \delta)$ -optimal  $t$ -policy の存在を与えたが, この § では更にそれが Markov  $t$ -policy の class の中にあることを示す。そのために, 次のように  $\{\beta_n^N(\pi)\}$  を  $\{v_n^N(\pi)\}$  に変換する。

$$\alpha^{n-1} v_n^N(\pi) = \beta_n^N(\pi) - \sum_{k=1}^{n-1} \alpha^{k-1} \gamma_k, \quad n=1, 2, \dots, N \quad (4.1)$$

従って  $\{\beta_n^N(\pi)\}$  の定義式 (3.1) は  $\{v_n^N(\pi)\}$  を用いて次のように書き換えられる。

$$\begin{cases} v_n^N(\pi) = \max \left[ g_n, \pi_n g_n \{ \alpha v_{n+1}^N(\pi) + \gamma_n \} \right], & n=1, 2, \dots, N-1 \\ v_N^N(\pi) = g_N \end{cases} \quad (4.2)$$

(注意) (4.2) をみると  $\{\beta_n^N(\pi)\}$ -process では、はっきりしなかった事が明らかになる。即ち (4.2) は policy  $\pi$  に従って行動してゆくとき、最適な stopping rule によって得られる optimal return に関する dynamic programming form の再帰式であり、 $v_n^N(\pi)$  は残りの stage 数  $n$  であるときの optimal return を表わす。

次にあげる Lemma は §3 で引用した Mackey の陰関数の定理に匹敵する重要なもう一つの可測陰関数の定理である。

Lemma 4.1 (Blackwell and Ryll-Nardzewski [1])

$X, Y$  を Borel set とする。任意の  $f \in Q(Y|X)$ , 任意の  $w \in M(XY)$ , 任意の  $\varepsilon > 0$  に対して

$$f[\{y; w(x, f(x)) \geq w(x, y) - \varepsilon\} | x] = 1 \quad \text{for } \forall x \in X$$

なる Baire 関数  $f: X \rightarrow Y$  が存在する。



次の Lemma もまた一つの陰関数の定理である。

Lemma 4.2 任意の  $\pi, p \in P(S)$ ,  $w \in M(SA)$ ,  $\epsilon > 0$   
 および各  $n$  ( $1 \leq n \leq N-1$ ) に対して

$$p\pi_1 g_1 \cdots \pi_{n-1} g_{n-1} \{f_n w \geq \pi_n w - \epsilon\} = 1$$

なる Baire 関数  $f_n; S \rightarrow A$  が存在する。

(略証)

$$\mathcal{H} \equiv p \in \pi \in P(SASA \cdots).$$

$\mathcal{H}$  の下で  $\Delta_n$  が与えられたときの条件のもとでの  $a_n$  についての conditional prob. measure を  $\pi_n^*$  とすると,

Lemma 4.1 により 任意の  $w \in M(SA)$ ,  $\epsilon > 0$  および各  $n$  に対して

$$\pi_n^* \{w(\Delta_n, f_n(\Delta_n)) \geq w(\Delta_n, a_n) - \epsilon \mid \Delta_n\} = 1$$

for  $\forall \Delta_n \in S$

なる Baire 関数  $f_n; S \rightarrow A$  が存在する。故に

$$\begin{aligned} \mathcal{H} \{w(\Delta_n, f_n(\Delta_n)) \geq w(\Delta_n, a_n) - \epsilon\} \\ = p\pi_1 g_1 \cdots \pi_{n-1} g_{n-1} \pi_n^* \{w(\Delta_n, f_n(\Delta_n)) \geq w(\Delta_n, a_n) - \epsilon \mid \Delta_n\} \\ = p\pi_1 g_1 \cdots \pi_{n-1} g_{n-1} \cdot 1 = 1 \end{aligned}$$

(Q.E.D.)

Theorem 4.1 任意の  $p \in P(S)$ ,  $\epsilon > 0$  および任意の  $t$ -policy  $(\pi, t)$  に対して,  $(\pi, t)$  は  $(p, \epsilon)$ -dominate する Markov

$\tau$ -policy が存在する。

(略証)

$(\pi, \tau)$  は任意の  $\tau$ -policy とする。

$$(\pi, \tau) \in \Lambda^N \text{ for some } N.$$

(4.1) から

$$\beta_1^N(\pi) = v_1^N(\pi).$$

Proposition 1, 2 により

$$E^\pi(x_\tau) \leq E^\pi(x_{\tau_N(\tau)}) = \beta_1^N(\pi) = v_1^N(\pi) \quad (4.3)$$

(4.2) より

$$\begin{aligned} v_{N-1}^N(\pi_{N-1}) &= \max [g_{N-1}, \pi_{N-1} \delta_{N-1} \{ \alpha g_N + \gamma_{N-1} \}] \\ &= \max [g_{N-1}, \pi_{N-1} w_{N-1}]. \end{aligned} \quad (4.4)$$

ここで

$$w_{N-1} = \delta_{N-1} \{ \alpha g_N + \gamma_{N-1} \} \in M(SA).$$

$$\zeta \equiv \epsilon / (1 + \alpha + \alpha^2 + \dots + \alpha^{N-1}) < \epsilon.$$

Lemma 4.2 により Baire 関数  $f_{N-1} : S \rightarrow A$  があって

$$\pi_{N-1} w_{N-1} \leq f_{N-1} w_{N-1} + \zeta \quad \text{with } p\pi_1 \delta_1 \dots \pi_{N-2} \delta_{N-2} \text{-prob. 1}$$

次に (4.4) から

$$\begin{aligned} v_{N-1}^N(\pi_{N-1}) &\leq \max [g_{N-1}, f_{N-1} w_{N-1} + \zeta] \\ &\quad \text{with } p\pi_1 \delta_1 \dots \pi_{N-2} \delta_{N-2} \text{-prob. 1} \\ &\leq \max [g_{N-1}, f_{N-1} w_{N-1}] + \zeta \\ &= v_{N-1}^N(f_{N-1}) + \zeta \end{aligned}$$

$$\begin{aligned} \therefore \pi_{N-2} g_{N-2} \{ \alpha v_{N-1}^N(\pi_{N-1}) + \gamma_{N-2} \} \\ \leq \pi_{N-2} w_{N-2} + \alpha \zeta \quad \text{with } p_{\pi_1 g_1 \dots \pi_{N-3} g_{N-3}} \text{-prob. } 1 \end{aligned}$$

ただし

$$w_{N-2} = g_{N-2} \{ \alpha v_{N-1}^N(f_{N-1}) + \gamma_{N-2} \} \in M(SA).$$

再び Lemma 4.2 により Barie 関数  $f_{N-2}; S \rightarrow A$  がとれ  
て,

$$\begin{aligned} v_{N-2}^N(\pi_{N-2}, \pi_{N-1}) &= \max [g_{N-2}, \pi_{N-2} g_{N-2} \{ \alpha v_{N-1}^N(\pi_{N-1}) + \gamma_{N-2} \}] \\ &\leq \max [g_{N-2}, f_{N-2} w_{N-2} + \zeta(1+\alpha)] \quad \text{with } p_{\pi_1 g_1 \dots \pi_{N-3} g_{N-3}} \text{-prob. } 1 \\ &\leq v_{N-2}^N(f_{N-2}, f_{N-1}) + \zeta(1+\alpha) \end{aligned}$$

以下同様に Barie 関数  $f_{N-3}, \dots, f_1$  をとらんと

$$v_1^N(\pi) \leq v_1^N(f_1, f_2, \dots, f_{N-1}) + \varepsilon \quad \text{with } p \text{-prob. } 1 \quad (4.5)$$

$$\pi^* \equiv \{ f_1, f_2, \dots, f_{N-1} \}$$

$$\tau^* \equiv \text{the first } n \text{ such that } v_n^N(\pi^*) = g_n$$

と置く

$$\tau^* = \tau_N(\pi^*). \quad (\text{ただし } \tau_N(\pi^*) \text{ は (3.2) で定義される})$$

かつ,  $(\pi^*, \tau^*)$  は Markov  $t$ -policy である。

(4.3) により (4.5) を書き直して

$$p \{ E^{\pi^*}(x_{\tau^*}) \geq v_1^N(\pi) - \varepsilon \} = 1$$

再び (4.3) により

$$p \{ E^{\pi^*}(x_{\tau^*}) \geq E^{\pi}(x_t) - \varepsilon \} = 1.$$

(Q.E.D.)

Theorem 4.2 Theorem 3.2 と同じ仮定をおく。

任意の  $p \in P(S)$ , 任意の  $\varepsilon > 0$ , 任意の  $\delta > 0$  に対して,  
 $(p, \varepsilon, \delta)$ -optimal Markov  $t$ -policy が存在する。

(略証)

Theorem 3.2 と Theorem 4.1 より 明らか。

Corollary 4.1 Theorem 3.2 と同じ仮定をおく。

(a)  $S$  が有限集合であれば, 任意の  $p \in P(S)$ ,  $\varepsilon > 0$  に対して,  
 $(p, \varepsilon)$ -optimal Markov  $t$ -policy が存在する。

(b)  $p \in P(S)$  が finite support をもてば, 任意の  $\varepsilon > 0$   
 に対して,  $(p, \varepsilon)$ -optimal Markov  $t$ -policy が存在する。

(注意) Theorem 4.1 の証明中で,  $v_n^N(\pi)$  と書くべき所も  
 $v_{N-1}^N(\pi_{N-1})$ ,  $\pi_{N-2}^N(\pi_{N-2}, \pi_{N-1})$  等と書いたが, 実は  $v_n^N(\pi)$   
 は  $\pi = \{\pi_1, \pi_2, \dots\}$  の成分の中,  $(\pi_n, \pi_{n+1}, \dots, \pi_{N-1})$   
 のみに depend するからである。

## §5 Stationary stopped-policy

この § では, D-case に対して, stationary  $s$ -policy の  
 存在と, それに関連して最適方程式 (optimality equation)

等について言論する。

この § ではつぎの仮定をおく。

- (i)  $0 \leq \alpha < 1$
- (ii)  $M(X)$  は  $X$  上のすべての有界 Baire 関数からなる空間を表わす。  $M(X)$  におけるノルムを  $\|u\| = \sup_{x \in X} |u(x)|$  で定義する。

(iii)  $g_1 = g_2 = g_3 = \dots = g \in Q(S|SA)$ .

次に degenerate  $f \in Q(A|S)$  に対して operator  $T_f$ ;  $M(S) \rightarrow M(S)$  を次式で定義する。更に  $A_f$  を定義する。

$$T_f u = fg(\gamma + \alpha u),$$

$$A_f u = \max [g, T_f u].$$

Markov policy  $\pi^N = \{f_1, f_2, \dots, f_N\}$  ( $N < \infty$ ) に対して operator  $U_{\pi^N}, L_{\pi^N}$ ;  $M(S) \rightarrow M(S)$  を次式で定義する。

$$U_{\pi^N} u = \max_{1 \leq n \leq N} T_{f_n} u,$$

$$L_{\pi^N} u = \max [g, U_{\pi^N} u].$$

これらの operator に對して直ちに 次の 2 つの Lemma を得る。

### Lemma 5.1

(a)  $T_f, A_f, U_{\pi^N}, L_{\pi^N}$  は monotone operator である。

(b)  $u \in M(S)$ , 定数  $c$  に対して

$$T_f(u+c) = T_f u + \alpha c, \quad U_{\pi^N}(u+c) = U_{\pi^N} u + \alpha c.$$

(c)  $u \in M(S)$ , 定数  $c > 0$  に対して

$$A_f(u+c) \leq A_f u + \alpha c, \quad L_{\pi^N}(u+c) \leq L_{\pi^N} u + \alpha c.$$

(d)  $u \in M(S)$ , 定数  $c < 0$  に対して

$$A_f(u+c) \geq A_f u + \alpha c, \quad L_{\pi^N}(u+c) \geq L_{\pi^N} u + \alpha c.$$

Lemma 5.2 (a)  $T_f, A_f, U_{\pi^N}, L_{\pi^N}$  は  $M(S)$  上の contraction mapping でその contraction coefficient は  $\alpha$  である.

(b)  $u, v \in M(S)$ , Markov policy  $\pi = \{f_1, f_2, \dots\}$  に対して

$$\lim_{n \rightarrow \infty} \|A_{f_1} A_{f_2} \cdots A_{f_n} u - A_{f_1} A_{f_2} \cdots A_{f_n} v\| = 0.$$

次に  $\pi^N$ -generated の定義を与える。

Markov policy  $\pi^N = \{f_1, f_2, \dots, f_N\}$  に対して

$f(S \rightarrow A)$  が  $\pi^N$ -generated ;  $S$  の Borel 分割  $S_1, S_2, \dots$  ;

$$S_N \text{ が有って } f = f_n \text{ on } S_n$$

Markov policy  $\hat{\pi}^N = \{g_1, g_2, \dots, g_N\}$  が  $\pi^N$ -generated ;

各  $g_n$  が  $\pi^N$ -generated

$F(\pi^N)$  ; すべての  $\pi^N$ -generated function の集合

$G^N(\pi^N)$  ;  $\Pi^N$  に属するすべての  $\pi^N$ -generated Markov policy の集合

Lemma 5.3

- (a)  $T_f u \leq U_{\pi^N} u$  for any  $u \in M(S)$ , for any  $f \in F(\pi^N)$   
 (b)  $A_f u \leq L_{\pi^N} u$  for any  $u \in M(S)$ , for any  $f \in F(\pi^N)$

Lemma 5.4 任意の  $u \in M(S)$  に対して次のような  $f \in F(\pi^N)$  が存在する:

$$T_f u = U_{\pi^N} u \quad \text{and} \quad A_f u = L_{\pi^N} u.$$

Lemma 5.5  $u_N^*$  を  $L_{\pi^N}$  の fixed point とする。

- (a) 任意の  $\varepsilon > 0$  に対して  $N$  を十分大きくとれば

$$\beta_1^{N+1}(\hat{\pi}^N) - \varepsilon \leq u_N^* \quad \text{for } \hat{\pi}^N \in G^N(\pi^N).$$

- (b) 任意の  $\varepsilon > 0$  に対して  $N$  を十分大きくとれば

$$\beta_1^{N+1}(f^{(N)}) \geq u_N^* - \varepsilon$$

ある  $f \in F(\pi^N)$  をえらぶことが出来る。ただし  $f^{(N)} = \{f, f, \dots, f\}$  ( $N$  factors)

(略証)

- (a)  $\hat{\pi}^N = \{\hat{f}_1, \hat{f}_2, \dots, \hat{f}_N\} \in G^N(\pi^N)$  とする。

Lemma 5.3 (b) より

$$A_{\hat{f}_i} u \leq L_{\pi^N} u \quad \text{for } 1 \leq i \leq N, \quad \text{for } u \in M(S)$$

$$\therefore A_{\hat{f}_N} u_N^* \leq L_{\pi^N} u_N^* = u_N^*$$

$$\therefore A_{\hat{f}_1} A_{\hat{f}_2} \cdots A_{\hat{f}_N} u_N^* \leq u_N^* \quad (5.1)$$

Lemma 5.2 (b) と (5.1) により,  $N$  を十分大きくとれば

$$\begin{aligned} u_N^* &\geq A_{\hat{f}_1} A_{\hat{f}_2} \cdots A_{\hat{f}_N} g - \varepsilon \\ &= v_1^{N+1}(\hat{\pi}^N) - \varepsilon = \beta_1^{N+1}(\hat{\pi}^N) - \varepsilon. \end{aligned}$$

(b) Lemma 5.4 (b) より  $f \in F(\pi^N)$  があって

$$A_f u_N^* = L_{\pi^N} u_N^* = u_N^*.$$

$$\therefore (A_f)^N u_N^* = u_N^*$$

Lemma 5.2 (b) により  $N$  を十分大きくとれば

$$(A_f)^N g + \varepsilon \geq (A_f)^N u_N^* = u_N^* \quad (5.2)$$

よるに  $(A_f)^N g = v_1^{N+1}(f^{(N)}) = \beta_1^{N+1}(f^{(N)})$  となるから (5.2) より,

$$\beta_1^{N+1}(f^{(N)}) \geq u_N^* - \varepsilon.$$

(Q.E.D.)

Theorem 5.1 (A3), (A4), (A5) を仮定する。このとき

任意の  $p \in P(S)$ , 任意の  $\varepsilon > 0$ , 任意の  $\delta > 0$  に対して  $(p, \varepsilon, \delta)$ -optimal stationary  $\delta$ -policy が  $\tau$ -policy として存在する。

(略証)

$(\pi^{N_1}, \tau_{N_1}) \in \Lambda^{N_1}$  を  $(p, \frac{\varepsilon}{3}, \delta)$ -optimal Markov  $\tau$ -policy とする。(存在は Theorem 4.2)

$\delta_3$  における議論から  $(\pi^{N_1}, \tau_{N_1})$  は  $(p, \frac{\varepsilon}{3})^{N_1}$ -optimal として得られる。

一方, 各  $N$  に対する  $(p, \frac{\varepsilon}{3})^N$ -optimal  $\tau$ -policy を



$(\pi^N, \tau_N)$  とすれば Lemma 3.5 より

$$\mathbb{P}\left\{E^{\pi^N}(x_{\tau_N}) \geq E^{\pi}(x_t) - \frac{\varepsilon}{3}\right\} \geq 1 - \delta \quad \text{for } N \geq N_1, \\ \text{for } (\pi, t) \in \Lambda \quad (5.3)$$

$L_{\pi^N}$  の fixed point を  $u_{N-1}^*$  とすれば Lemma 5.5 (b) より

$N_2$  を十分大きくとれば, 各  $N \geq N_2$  に対して

$$\exists f \in F(\pi^N); \beta_1^N(f^{(N-1)}) \geq u_{N-1}^* - \frac{\varepsilon}{3} \quad (5.4)$$

Lemma 5.5 (a) より  $N_3$  を十分大きくとれば, 各  $N \geq N_3$  に対して

$$u_{N-1}^* \geq \beta_1^N(\pi^N) - \frac{\varepsilon}{3} \quad (5.5)$$

$N_0 = \max[N_1, N_2, N_3]$  とおけば (5.4) (5.5) より

$$\beta_1^{N_0}(f^{(N_0-1)}) \geq u_{N_0-1}^* - \frac{\varepsilon}{3} \geq \beta_1^{N_0}(\pi^{N_0}) - \frac{2}{3}\varepsilon \quad (5.6)$$

そこで

$$\tilde{\tau}_{N_0} = \text{the first } n \text{ such that } \beta_n^{N_0}(f^{(N_0-1)}) = x_n$$

とおけば

$$\tilde{\tau}_{N_0} \in C^{N_0}(f^{(N_0-1)}) \Rightarrow \beta_1^{N_0}(f^{(N_0-1)}) = E^{f^{(N_0-1)}}(x_{\tilde{\tau}_{N_0}})$$

さらに (5.6) より

$$E^{f^{(N_0-1)}}(x_{\tilde{\tau}_{N_0}}) \geq E^{\pi^{N_0}}(x_{\tau_{N_0}}) - \frac{2\varepsilon}{3} \quad (5.7)$$

(5.3) (5.7) より

$$\mathbb{P}\left\{E^{f^{(N_0-1)}}(x_{\tilde{\tau}_{N_0}}) \geq E^{\pi}(x_t) - \varepsilon\right\} \geq 1 - \delta \quad \text{for } \forall (\pi, t) \in \Lambda$$

即ち  $(f^{(N_0-1)}, \tilde{\tau}_{N_0})$  は  $(\rho, \varepsilon, \delta)$ -optimal stationary  $t$ -policy

である。

(Q.E.D.)

Corollary 5.1 Theorem 5.1 と同じ仮定をおく。

- (a)  $S$  が有限集合なら, 任意の  $p \in P(S)$ ,  $\varepsilon > 0$  に対して  $(p, \varepsilon)$ -optimal stationary  $t$ -policy が存在する。
- (b)  $p \in P(S)$  が finite support をもてば, 任意の  $\varepsilon > 0$  に対して,  $(p, \varepsilon)$ -optimal stationary  $t$ -policy が存在する。

次に Markov policy  $\pi = \{f_1, f_2, \dots\}$  に対して operator  $U_\pi, L_\pi$  を定義する。即ち

$$U_\pi u = \sup_n T_{f_n} u, \quad u \in M(S)$$

$$L_\pi u = \max[g, U_\pi u].$$

このとき,  $U_\pi, L_\pi$  は monotone operator であり,  $M(S)$  上の contraction mapping になる。

Markov policy  $\pi = \{f_1, f_2, \dots\}$  に対して

$f; S \rightarrow A$  が  $\pi$ -generated;  $S$  の Borel partition  $S_1, S_2, \dots$

が存在して  $f = f_n$  on  $S_n$ .

Markov policy  $\hat{\pi} = \{g_1, g_2, \dots\}$  が  $\pi$ -generated;

$\sum g_n$  が  $\pi$ -generated

$F(\pi)$ ; すべての  $\pi$ -generated function の集合

$G(\pi)$ ;  $\Pi$  に属するすべての  $\pi$ -generated M-policy の集合

$\Lambda_\pi \equiv \{(\hat{\pi}, \hat{c}) \mid \hat{\pi} \in G(\pi), \hat{c} \in C(\hat{\pi})\}$

## Optimality equation

$A_a$  を  $f \equiv a$  なる  $f$  に対する operator とする。この  
とき  $u \in M(S)$  に対して

$$u = \sup_{a \in A} A_a u \quad (5.8)$$

が成立するとき、 $u$  は optimality equation をみたすといふ、  
(5.8) を optimality equation といふ。

## Optimal return

$u \in M(S)$  が

$$(i) \quad u \geq E^{\hat{\pi}}(x_{\hat{e}}) \quad \text{for } \forall (\hat{\pi}, \hat{e}) \in \Lambda_{\pi}$$

$$(ii) \quad \forall \varepsilon > 0, \exists (\hat{\pi}, \hat{e}) \in \Lambda_{\pi} ; E^{\hat{\pi}}(x_{\hat{e}}) \geq u - \varepsilon$$

をみたすとき、 $u$  は optimal return in  $\Lambda_{\pi}$  といふ。

$u \in M(S)$  が

$$(iii) \quad u \geq E^{\hat{\pi}}(x_{\hat{e}}) \quad \text{for } \forall (\hat{\pi}, \hat{e}) \in \Lambda$$

$$(iv) \quad \forall \varepsilon > 0, \exists (\hat{\pi}, \hat{e}) \in \Lambda ; E^{\hat{\pi}}(x_{\hat{e}}) \geq u - \varepsilon$$

をみたすとき、 $u$  は optimal return といふ。

Theorem 5.2 (A3), (A4), (A5) をおく。

(a)  $\pi$  を任意の policy,  $L_{\pi}$  の fixed point を  $u^*$  とすれば  
 $u^*$  は optimal return in  $\Lambda_{\pi}$  である。

(b)  $\varepsilon \geq 0$  とする。もし  $\varepsilon$ -optimal  $s$ -policy があれば

各  $\varepsilon' > 0$  に対して  $(\varepsilon/(1-\alpha) + \varepsilon')$ -optimal stationary  $s$ -policy が存在する。

(c)  $u \in M(S)$  が  $A_a u \leq u$  for all  $a \in A$  をみたせば  $E^\pi(x_t) \leq u$  for  $(\pi, t) \in \Lambda$  が成立する。

(d) 各  $\varepsilon > 0$  に対して  $\varepsilon$ -optimal  $s$ -policy があれば optimal return は optimality equation をみたす。

(証明は非常に長いので省略)

$\pi$  が semi-Markov ;  $\pi = \{\pi_1, \pi_2, \dots\}$  において  $\pi_n \in \mathcal{Q}(A|SS)$  かつ  $\pi_n$  は degenerate

semi-Markov policy を用いて, 次の定理が導かれる。

Theorem 5.3 (A3) (A4) (A5) をおく。

$s$ -policy  $(\pi^*, t^*)$  が optimal

$\iff E^{\pi^*}(x_{t^*})$  が optimality equation をみたす

Theorem 5.4 optimal equation は高々一つの bounded solution をもつ。

## §6 Additional results.

特に, action space が countable set または finite set のときを扱う。

Theorem 6.1 (A3) (A4) (A5) をおく。

(a)  $A$  が countable set である, 任意の  $\varepsilon > 0$  に対して  $\varepsilon$ -optimal stationary  $\pi$ -policy が存在する。

(b) (A3) (A5) の他に ((A4) は不要), (i)  $\lim_{n \rightarrow \infty} x_n'' \equiv \lim_{n \rightarrow \infty} \sum_{k=1}^n \alpha^{k-1} y_k'' = +\infty$ , (ii)  $\{x_n^-\}$  が  $e_\pi$  について一様可積分, (iii)  $\sup_N E^{f^\infty}(x'_{\tau_N(f^\infty)}) < \infty$  for  $\forall f^\infty$  の仮定をおく。このとき  $A$  が finite set なら stationary stopping time をもつような optimal stationary  $s$ -policy が存在する。

### [付記]

この報告は [2] からの抜粋で, ここに紹介しなかったかなりの部分, および証明を除いた箇所については [2] を見て頂きたい。

なお §6 は optimal 又は  $\varepsilon$ -optimal  $s$ -policy の存在定理であるが, action space が countable でなくとも, 一般に compact である, 他に transition probability  $p$  に弱連続の仮定等を加えれば Theorem 6.1 の結果と同様の結果が得られることが, 最近 S. Iwamoto によって証明された。([3])

## REFERENCES

- [1] D. Blackwell and C. Ryll-Nardzewski ; Non-existence of everywhere proper conditional distributions. Ann. Math. Statist. 34 (1963), 223-225.
- [2] N. Furukawa and S. Iwamoto ; Stopped decision processes on complete separable metric spaces.  
(to appear)
- [3] S. Iwamoto ; Stopped decision processes on compact metric spaces. (to appear)
- [4] H.J. Kushner ; Computational procedures for optimal stopping problems for Markov chains. Jour. Math. Anal. Appl. 25 (1969), 607-615.
- [5] G.W. Mackey ; Borel structure in groups and their duals. Trans. Amer. Math. Soc. 85 (1957), 134-165.
- [6] K.R. Parthasarathy ; Probability measures on metric spaces. (1967), Academic Press.